
A Multi-Token Coordinate Descent Method for Vertical Federated Learning

Pedro Valdeira^{1,2}, Yuejie Chi¹, Cláudia Soares³, João Xavier²

¹Carnegie Mellon University, ²Instituto Superior Técnico, ³NOVA School of Science and Technology
pvaldeira@cmu.edu

Abstract

Communication efficiency is a major challenge in federated learning. In client-server schemes, the server constitutes a bottleneck, and while decentralized setups spread communications, they do not reduce them. We propose a communication-efficient algorithm for semi-decentralized vertical federated learning dealing with feature-distributed data. Our multi-token method can be seen as a parallel Markov chain (block) coordinate descent algorithm. In this work, we formalize the multi-token semi-decentralized scheme, which subsumes the client-server and decentralized setups, and design a feature-distributed learning algorithm for this setup. Numerical results show the improved communication efficiency of our algorithm.

1 Introduction

Federated Learning (FL) is a machine learning paradigm where data is distributed across a set of clients who collaborate to learn a model without sharing local data [23]. Most FL literature considers data distributed by samples (horizontal FL), where each client holds all the features of a subset of the samples, yet recently there has been a growing interest on feature-distributed setups (vertical FL), where each client holds a subset of the features for all samples [1, 5, 7, 13].

FL often deals with the client-server setup with a star-shaped topology. However, such schemes have a single point of failure and suffer from a communication bottleneck on the server [19]. On the other hand, there is extensive literature on decentralized optimization—from earlier work motivated by applications such as wireless sensor networks and multiagent control [9, 25, 28], to recent work motivated by FL [17, 18]. Yet, these algorithms converge slowly in sparse and large networks [26] and, although they spread the communication load across the network, they tend to have a poor communication efficiency [37].

When concerned with the communications between clients, the use of a token method [3, 14, 16, 22, 24, 29], where a model-describing *token* follows a random walk over a communication graph (undergoing local updates), allows for better communication efficiency [14] than the more common consensus-based methods [9, 17, 25, 28], where asymptotic agreement is reached through successive local averaging. Yet, the convergence rate of token methods degrades even faster for larger and sparser networks, due to a lack of parallel communications. Multi-token methods [6, 14, 34] mitigate this problem by running multiple tokens simultaneously and combining them.

Motivated by the above observations, we propose a Semi-Decentralized FL (SDFL) multi-token algorithm for vertical FL. By using both client-server *and* client-client communications, we reduce the communications at the server while mitigating the slow convergence of decentralized algorithms in sparse and large networks. Our algorithm might be of particular interest for applications, for example, using time series data measured by personal devices to learn a model of some “cross-client”

phenomenon of interest (e.g. meteorological). Here, each sample would correspond to the data collected across the devices at a given timestamp.

Our main contributions are as follows.

- We formalize the multi-token semi-decentralized federated learning scheme, which is robust to server failures and flexible in the degree of dependence on the server, recovering both client-server and decentralized FL as particular cases.
- We further design a Multi-Token Coordinate Descent (MTCD) algorithm for vertical FL with improved communication efficiency over state-of-the-art methods.
- Numerical experiments are performed on both synthetic and real data for a variety of communication setups, showing the improved communication efficiency of MTCD as well as the effectiveness of resorting to multiple tokens.

Related works. Recently, SDFL approaches have been proposed to lower communication costs and deal with data heterogeneity [11, 20] and to handle intermittent connections, latency, and stragglers [4, 35]. Some SDFL works deal with (multi-layered) hierarchical networks [15, 36].

Coordinate Descent (CD) methods [33], where (blocks of) coordinates are updated sequentially, rather than simultaneously, are natural candidates for optimization in feature-distributed learning. The block selection is most often cyclic [2] or independent and identically distributed random [27, 30]. In contrast, [32] considers block selection following a Markov chain. Several extensions to CD have been proposed, such as acceleration and parallelization [10] and distributed CD methods [6, 21].

2 Single-token coordinate descent (STCD)

We now introduce a simple, particular case of our algorithm, closely related to [32] and the application mentioned therein, taken from [22]. Yet, we work in the primal domain and on a feature-distributed setting. In this section, we use the terms client and agent interchangeably.

2.1 Problem statement

Model. Our goal is to minimize the following regularized generalized linear model:

$$f(\boldsymbol{\theta}) \triangleq \sum_{n=1}^N \ell(\mathbf{x}_n^\top \boldsymbol{\theta}) + r(\boldsymbol{\theta}), \quad (1)$$

where $\boldsymbol{\theta}$ is a parameter of interest and we have N samples \mathbf{x}_n , $n \in \{1, \dots, N\}$. Our loss function ℓ is a (possibly nonconvex) smooth and block smooth loss function with a nonempty set of minimizers and r is a separable (possibly nonsmooth) closed proper function, $r(\boldsymbol{\theta}) = \sum_{k=1}^K r_k(\boldsymbol{\theta}_k)$, which we assume to have a simple proximal mapping.

Setup. In this section, we do not require the existence of a server. The clients $[K] \triangleq \{1, \dots, K\}$ learn the model in a fully decentralized manner, communicating through channels described by a static, undirected communication graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = [K]$ is the vertex set and \mathcal{E} the edge set. We assume \mathcal{G} is connected and denote the set of neighbors of agent k by $\mathcal{N}_k \triangleq \{i: \{i, k\} \in \mathcal{E}\}$.

Our dataset $\mathbf{X} \in \mathbb{R}^{N \times d}$ with samples $\mathbf{x}_n \in \mathbb{R}^d$, $n \in [N]$, is distributed by features across the clients. More precisely, each sample \mathbf{x}_n is partitioned as $\mathbf{x}_n = (\mathbf{x}_{nk})_{k=1}^K \in \mathbb{R}^d$, where $\mathbf{x}_{nk} \in \mathbb{R}^{d_k}$ and $\sum_{k=1}^K d_k = d$. Consequently, we can equivalently define the partitioning as $\mathbf{X} = [\mathbf{X}_1 \dots \mathbf{X}_K]$, where $\mathbf{X}_k \in \mathbb{R}^{N \times d_k}$. Our parameter $\boldsymbol{\theta} \in \mathbb{R}^d$ is partitioned similarly to \mathbf{x}_n , that is, $\boldsymbol{\theta} = (\boldsymbol{\theta}_k)_{k=1}^K \in \mathbb{R}^d$ with $\boldsymbol{\theta}_k \in \mathbb{R}^{d_k}$. The labels $\mathbf{y} \in \mathbb{R}^N$ are part of the loss function, which we assume all agents to know.

2.2 Algorithm

We now propose a decentralized token method algorithm to minimize (1). Our *token* carries $\mathbf{z} \triangleq \mathbf{X}\boldsymbol{\theta} \in \mathbb{R}^N$ and performs a random walk over the graph, being updated and communicated by each agent after performing local computations. We denote by k^t the agent holding the token at iteration t .

Algorithm 1 Single-token Coordinate Descent (STCD)

Initialize $\theta^0 = \mathbf{0}$, $\mathbf{z}^0 = \mathbf{0}$, arbitrary k^0 , and choose step-size μ
for $t \geq 0$ **do**
 if $t \bmod Q = 0$ **then**
 Agent k^t sends \mathbf{z}^{t+1} to agent $k^{t+1} \sim \mathcal{U}(\mathcal{N}_{k^t})$
 else
 $k^{t+1} = k^t$
 end if
 $\theta_k^{t+1} = \begin{cases} \text{prox}_{\mu r_k}(\theta_k^t - \mu \nabla_k L(\theta^t)) & \text{if } k = k^t \\ \theta_k^t & \text{if } k \neq k^t \end{cases} \quad \{\mathbf{z}^t \text{ used when computing } \nabla_{k^t} L(\theta^t)\}$
 $\mathbf{z}^{t+1} = \mathbf{z}^t + \mathbf{X}_{k^t}(\theta_{k^t}^{t+1} - \theta_{k^t}^t)$
end for

The token suffices to compute partial gradients. From the definition of \mathbf{z} , we have at iteration t that $\mathbf{z}^t = \sum_{k=1}^K \mathbf{X}_k \theta_k^t$. Letting $L(\theta) \triangleq \sum_{n=1}^N \ell(\mathbf{x}_n^\top \theta)$, we see that, given \mathbf{z}^t , agent k can compute the gradient of the loss function with respect to (w.r.t.) its local parameter θ_k :

$$\nabla_k L(\theta) \triangleq \nabla_{\theta_k} \sum_{n=1}^N \ell(\mathbf{x}_n^\top \theta) = \sum_{n=1}^N \nabla_z \ell(z_n) \mathbf{x}_{nk}, \quad z_n = \mathbf{x}_n^\top \theta,$$

allowing agent k to (locally) perform a proximal CD step w.r.t. θ_k , as long as it also holds θ_k^t :

$$\theta_k^{t+1} = \text{prox}_{\mu r_k}(\theta_k^t - \mu \nabla_k L(\theta^t)), \quad \text{prox}_g(\cdot) \triangleq \text{argmin}_u g(u) + \frac{1}{2} \|u - \cdot\|^2.$$

Given that we only update θ_k at agent k , k^t always knows $\theta_{k^t}^t$.

We now describe STCD, summarized in Algorithm 1, where \mathcal{U} denotes the uniform distribution.

Initialization. To start STCD at a given client k^0 , we need to ensure that \mathbf{z}^0 which, as we saw, contains all the information needed to do a (proximal) CD step, and θ^0 verify $\mathbf{z}^0 = \mathbf{X} \theta^0$. To do this without requiring any information from the other clients, we must initialize $\theta^0 = \mathbf{0}$ and $\mathbf{z}^0 = \mathbf{0}$.

Updating and communicating the token. After performing a local (proximal) coordinate descent step, which gives agent k^t the iterate $\theta_{k^t}^{t+1}$, we need to update the token to ensure that $\mathbf{z}^{t+1} = \mathbf{X} \theta^{t+1}$. Observing that \mathbf{z}^t allows for a local update $\mathbf{z}^{t+1} = \mathbf{z}^t + \mathbf{X}_{k^t}(\theta_{k^t}^{t+1} - \theta_{k^t}^t)$ whose result can then be sent to agent k^{t+1} , we see that, by induction, \mathbf{z}^t can be kept up-to-date throughout our algorithm. Additionally, this also means that an agent can actually perform Q multiple coordinate descent steps w.r.t. its local parameters, θ_k .

In essence, STCD is a technique allowing for Markov Chain Coordinate Descent [32] to be performed in feature-distributed setups. In terms of the progress made in the parameter space, Algorithm 1 differs from [32] only in that it requires initializing $\theta^0 = \mathbf{0}$ and in that we allow for $Q > 1$.

Let f be a function that is L_g -smooth and block L -smooth for all blocks k , with a nonempty set of minimizers, and let r be a separable closed proper function. It is shown in [32] that, for $Q = 1$, Algorithm 1 converges almost surely to a solution of (1).

2.3 Limitations

The decentralized token algorithm in this section has an appealing simplicity. However, while it outperforms state-of-the-art feature-distributed learning algorithms in a variety of setups, as we will see in Section 4, its performance deteriorate faster with network connectivity than these decentralized consensus-based algorithms. Yet, this simple algorithm acts as a stepping stone to the more general multi-token algorithm we now present in Section 3, which mitigates this problem.

3 Multi-token coordinate descent (MTCDD)

In this section, instead of a fully decentralized setup, we deal with SDFL and generalize the algorithm used in the previous section in a direction allowing us to leverage the existence of a server.

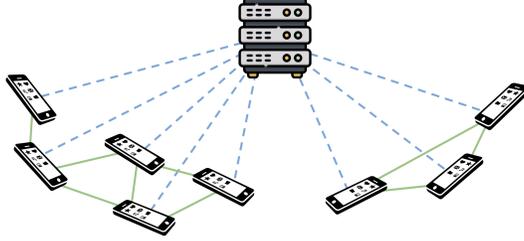


Figure 1: Semi-decentralized setup with $K = 9$. Client-server communications are represented by the dashed blue lines and client-client communications are represented by the unbroken green lines.

3.1 Problem statement

In this section, the set of clients and the data and model partitioning are similar to Section 2. However, we now consider a server-enhanced setup (that is, SDFL), illustrated in Figure 1, where we also take into account a central server with links to all the clients. In this section, we do not require the graph describing the communications between the clients to be connected.

Considering the same model (1) as in the previous section, we now develop an algorithm to leverage the existence of this server, exploiting it to mitigate the deterioration of the performance of Algorithm 1 in poorly connected networks.

3.2 Algorithm

The algorithm introduced in this section can be seen as having two parts. The roaming part is similar to the procedure in Algorithm 1, except for the fact that we run multiple tokens in parallel, each with an associated model. The novelty lies in the averaging part, where we exploit the access to the server by periodically averaging these tokens there.

More precisely, we use Γ tokens $z_\gamma \in \mathbb{R}^N$, $\gamma \in [\Gamma]$, each with an associated model instance θ_γ and $z_\gamma \triangleq \mathbf{X}\theta_\gamma \in \mathbb{R}^N$. Similarly to what we had in Algorithm 1, we initialize $\theta_\gamma^0 = \mathbf{0}$ and $z_\gamma^0 = \mathbf{0}$, $\gamma \in [\Gamma]$, and for $S \times Q$ iterations, each such token undergoes a process similar to the one in Algorithm 1. The difference lies in the fact that now, every S hops of the tokens, between each of which each agent performs Q local steps, the Γ tokens are averaged at the server and the Γ models are averaged locally, at the clients. Crucially, these averaging operations allow us to preserve the relationship defining \mathbf{Z}_γ^t . (This relationship is also preserved during the roaming part, as explained in the previous section.)

This method, a parallel Markov Chain Coordinate Descent for feature-distributed SDFL setups, is summarized in Algorithm 2, where k_γ^t denotes the agent holding token γ at time t .

Recovering client-server and decentralized setups. If no client-server communications are available ($S \rightarrow \infty$) our algorithm is reduced to a decentralized one. In this setting, even if at a lower rate, our algorithm will still converge, as long as \mathcal{G} is connected, being reduced to Γ simultaneous runs of Algorithm 1. While a decentralized asynchronous token averaging is possible (by exploring the random intersections of multiple tokens are at the same client), this is significantly less effective than server averaging in terms of mitigating the effect of poor connectivity. In contrast, if the set of edges \mathcal{E} is empty and we assign a token per agent, we recover the client-server setting.¹

Extension to locally nonlinear models. We can generalize our model while keeping an additive structure w.r.t. $\{\mathbf{X}_k, \theta_k\}$, capturing nonlinearities between \mathbf{X}_k and θ_k . Letting $\mathbf{h}_k: \mathbb{R}^{|\theta_k| \times d_k} \rightarrow \mathbb{R}^E$, we would minimize:

$$f(\theta) \triangleq \sum_{n=1}^N \ell \left(\sum_{k=1}^K \mathbf{h}_k(\theta_k, \mathbf{x}_{nk}) \right) + r(\theta),$$

which includes generalized linear models as the particular case $\mathbf{h}_k(\theta_k, \mathbf{x}_{nk}) = \langle \theta_k, \mathbf{x}_{nk} \rangle$, where $E = 1$ and $|\theta_k| = |\mathbf{x}_{nk}| = d_k$. However, this would require additional communications, since

¹To be precise, we would have to change $k_\gamma^{t+1} \sim \mathcal{U}(\mathcal{N}_{k_\gamma^t})$ to $k_\gamma^{t+1} \sim \mathcal{U}(\bar{\mathcal{N}}_{k_\gamma^t})$, where $\bar{\mathcal{N}}_{k_\gamma^t} \triangleq \mathcal{N}_{k_\gamma^t} \cup \{k_\gamma^t\}$.

Algorithm 2 Multi-token Coordinate Descent (MTCD)

Initialize $\theta_\gamma^0 = \mathbf{0}$, $z_\gamma^0 = \mathbf{0}$, and arbitrary k_γ^0 , for $\gamma \in [\Gamma]$, and choose step-size μ

for $t \geq 0$ **do**

if $t \bmod Q = 0$ and $\lfloor t/Q \rfloor \bmod S = 0$ **then**

 clients send z_γ^t to the server

 server computes $z_1^t, \dots, z_\Gamma^t = \frac{1}{\Gamma} \sum_{\gamma=1}^\Gamma z_\gamma^t$

 server sends updated z_γ^t to the clients

for $k \in [K]$ **in parallel do**

$\theta_{1k}, \dots, \theta_{\Gamma k} = \frac{1}{\Gamma} \sum_{\gamma=1}^\Gamma \theta_{\gamma k}$

end for

end if

if $t \bmod Q = 0$ **then**

for $\gamma \in [\Gamma]$ **in parallel do**

 Agent k_γ^t sends z_γ^{t+1} to agent $k_\gamma^{t+1} \sim \mathcal{U}(\mathcal{N}_{k_\gamma^t})$

end for

else

$k_\gamma^{t+1} = k_\gamma^t, \gamma \in [\Gamma]$

end if

for $\gamma \in [\Gamma]$ **in parallel do**

$\theta_{\gamma k}^{t+1} = \begin{cases} \text{prox}_{\mu r_k}(\theta_{\gamma k}^t - \mu \nabla_k L(\theta_\gamma^t)) & \text{if } k = k_\gamma^t \\ \theta_{\gamma k}^t & \text{if } k \neq k_\gamma^t \end{cases} \{z_\gamma^t \text{ used when computing } \nabla_{k_\gamma^t} L(\theta_\gamma^t)\}$

$z_\gamma^{t+1} = z_\gamma^t + \mathbf{X}_k(\theta_{\gamma k}^{t+1} - \theta_{\gamma k}^t), k = k_\gamma^t$

end for

end for

a separate averaging of θ^t and z^t would no longer preserve their relation, leading to a need to recompute the tokens, instead of simply updating them in an online manner. This would require all clients to communicate with the server.

4 Experiments

In this section, we compare STCD to Dual Consensus Proximal Algorithm (DCPA) [1], a state-of-the-art decentralized method used as a baseline, and to MTCD.

Models and datasets. We perform ridge regression on a dataset generated following the same process as [1],² where the number of samples and the dimensionality are $N = 1000$ and $d = 2000$, respectively. We have the following objective function, with $\alpha = 10$:

$$f(\theta) = \frac{1}{2} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \frac{\alpha}{2} \|\theta\|_2^2.$$

We also perform sparse logistic regression on the Gisette dataset [12], where $N = 6000$ and $d = 5000$. Letting $s(z) \triangleq (1 + e^{-z})^{-1}$, we consider the following objective function, with $\beta = 1$:

$$f(\theta) = - \sum_{n=1}^N [y_n \log s(\mathbf{x}_n^\top \theta) + (1 - y_n) \log(1 - s(\mathbf{x}_n^\top \theta))] + \beta \|\theta\|_1, \quad y_n \in \{0, 1\}.$$

Metrics. We use CVXPY [8] to obtain f^* , and then compute the suboptimality gap $f(\theta^t) - f^*$ throughout our experiments. In the previous sections, we considered each local CD step as an STCD iteration, to avoid the visual clutter of having a double counter (for hops and local CD updates). In this section, to allow for a fairer comparison with DCPA, where only communication rounds count as a round (despite the fact that we need to compute a proximal operator with a gradient descent

²Each entry of \mathbf{X} is drawn uniformly at random from $\{0, 1\}$ and each entry of \mathbf{y} is drawn from a standard normal distribution (in both cases the entries are independent).

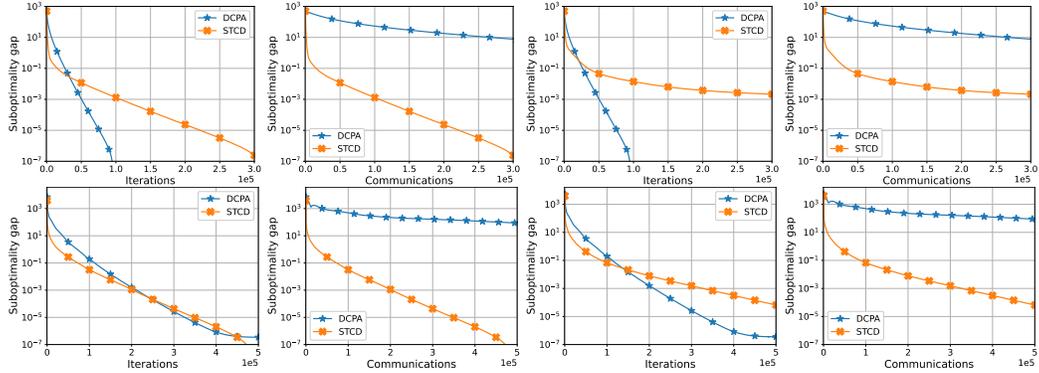


Figure 2: Suboptimality with respect to iterations and communications for a random network with $p = 0.4$ (four figures on the left) and a chain network (four figures on the right), both with $K = 20$. The top row concerns ridge regression and the bottom row concerns sparse logistic regression. Each communication is size N : STCD has one communication per iteration and DCPA has $2|\mathcal{E}|$.

subroutine), we count as an STCD iteration each hop in the network. We use the same definition of iteration for MTC.

Comparison with state-of-the-art. For ridge regression, we randomly generate a graph as in [31], with connectivity ratio $p = 0.4$, and a chain graph, both with $K = 20$ nodes and $d_k = 100$ for all k . For STCD, we use $\mu = 10^{-5}$ and $Q = 20$. For DCPA, we use $\mu_w = 0.01, \mu_y = 0.0003, \mu_x = 0.03$.

For sparse logistic regression, we similarly generate a random graph and chain graph, both with $K = 20$ nodes (now $d_k = 250$ for all k). For STCD, we use $\mu = 10^{-4}$ and $Q = 30$. For DCPA, we use $\mu_w = 0.001, \mu_y = 0.00003, \mu_x = 0.003$.

In Figure 2, we see that, while STCD does not improve upon DCPA in terms of progress per iteration, it significantly outperforms it in terms of communication efficiency. Yet, we can also see that STCD is particularly vulnerable to poorly connected networks, as evidenced by a performance deterioration when going from a random network to a chain network.

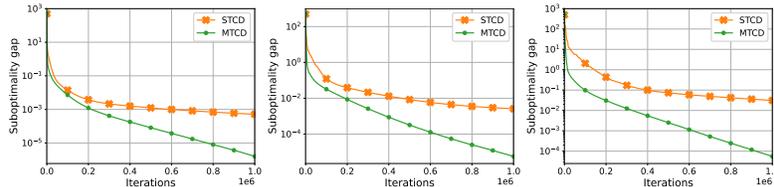


Figure 3: Suboptimality with respect to iterations for $K = 20, 40, 80$ chain networks (K increasing from left to right). All three plots concern the ridge regression model.

Comparing STCD with MTC. We again deal with the ridge regression problem above and resort to the same dataset. However, we now focus on chain graphs (for MTC, complemented by client-server communications), as they are poorly connected setups where the drawback of STCD is evident. We use $\mu = 10^{-5}$ and $Q = 20$ for both STCD and MTC and use $\Gamma = 10, S = 10$ for the latter throughout all three experiments, varying only the number of nodes $K \in \{20, 40, 80\}$.

In Figure 3, we see that, while both algorithms see a drop in performance for graphs with a lower connectivity, MTC mitigates this effect, consistently outperforming STCD.

5 Conclusions

We formalize the multi-token SDFL scheme and propose a communication-efficient SDFL algorithm for feature-distributed data. Numerical results show the improved communication efficiency of our algorithm as well as the power of endowing decentralized methods with periodical client-server communications. A natural extension to this work is a detailed study of how the number of tokens and the frequency of their averaging influence the convergence rate, both empirically and analytically.

Acknowledgments and Disclosure of Funding

This work is supported in part by the Fundação para a Ciência e a Tecnologia through the Carnegie Mellon Portugal Program; by the grant U.S. National Science Foundation CCF-2007911; and by the CMU-Portugal project CMU/TIC/0016/2021.

References

- [1] S. A. Alghunaim, Q. Lyu, M. Yan, and A. H. Sayed. Dual consensus proximal algorithm for multi-agent sharing problems. *IEEE Transactions on Signal Processing*, 69:5568–5579, 2021.
- [2] A. Beck and L. Tetruashvili. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- [3] D. P. Bertsekas. A new class of incremental gradient methods for least squares problems. *SIAM Journal on Optimization*, 7(4):913–926, 1997.
- [4] J. Bian and J. Xu. Mobility improves the convergence of asynchronous federated learning. *arXiv:2206.04742*, 2022.
- [5] T. J. Castiglia, A. Das, S. Wang, and S. Patterson. Compressed-VFL: Communication-efficient learning with vertically partitioned data. In *International Conference on Machine Learning*, pages 2738–2766. PMLR, 2022.
- [6] H. Chen, Y. Ye, M. Xiao, and M. Skoglund. Asynchronous parallel incremental block-coordinate descent for decentralized machine learning. *arXiv:2202.03263*, 2022.
- [7] T. Chen, X. Jin, Y. Sun, and W. Yin. VAFL: a method of vertical asynchronous federated learning. *arXiv:2007.06081*, 2020.
- [8] S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [9] J. C. Duchi, A. Agarwal, and M. J. Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic Control*, 57(3):592–606, 2012.
- [10] O. Fercoq and P. Richtárik. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on Optimization*, 25(4), 2015.
- [11] Y. Guo, Y. Sun, R. Hu, and Y. Gong. Hybrid local SGD for federated learning with heterogeneous communications. In *International Conference on Learning Representations*, 2021.
- [12] I. Guyon, S. Gunn, A. Ben-Hur, and G. Dror. Result analysis of the NIPS 2003 feature selection challenge. *Advances in neural information processing systems*, 17, 2004.
- [13] L. He, A. Bian, and M. Jaggi. COLA: Decentralized linear learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- [14] H. Hendrikx. A principled framework for the design and analysis of token algorithms. *arXiv:2205.15015*, 2022.
- [15] S. Hosseinalipour, S. S. Azam, C. G. Brinton, N. Michelusi, V. Aggarwal, D. J. Love, and H. Dai. Multi-stage hybrid federated learning over large-scale D2D-enabled fog networks. *IEEE/ACM Transactions on Networking*, 30(4):1569–1584, 2022.
- [16] B. Johansson, M. Rabi, and M. Johansson. A randomized incremental subgradient method for distributed optimization in networked systems. *SIAM Journal on Optimization*, 20(3):1157–1170, 2010.
- [17] A. Koloskova, N. Loizou, S. Boreiri, M. Jaggi, and S. Stich. A unified theory of decentralized SGD with changing topology and local updates. In *International Conference on Machine Learning*, pages 5381–5393. PMLR, 2020.
- [18] B. Li, S. Cen, Y. Chen, and Y. Chi. Communication-efficient distributed optimization in networks with gradient tracking and variance reduction. *Journal of Machine Learning Research*, 21:1–51, 2020.
- [19] X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu. Can decentralized algorithms outperform centralized algorithms? A case study for decentralized parallel stochastic gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017.

- [20] F. P.-C. Lin, S. Hosseinalipour, S. S. Azam, C. G. Brinton, and N. Michelusi. Semi-decentralized federated learning with cooperative D2D local model aggregations. *IEEE Journal on Selected Areas in Communications*, 39(12):3851–3869, 2021.
- [21] Y. Liu, X. Zhang, Y. Kang, L. Li, T. Chen, M. Hong, and Q. Yang. FedBCD: A communication-efficient collaborative learning framework for distributed features. *IEEE Transactions on Signal Processing*, 70:4277–4290, 2022.
- [22] X. Mao, K. Yuan, Y. Hu, Y. Gu, A. H. Sayed, and W. Yin. Walkman: A communication-efficient random-walk algorithm for decentralized optimization. *IEEE Transactions on Signal Processing*, 68:2513–2528, 2020.
- [23] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [24] A. Nedic and D. P. Bertsekas. Incremental subgradient methods for nondifferentiable optimization. *SIAM Journal on Optimization*, 12(1):109–138, 2001.
- [25] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [26] A. Nedic, A. Olshevsky, and M. G. Rabbat. Network topology and communication-computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976, 2018.
- [27] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [28] G. Qu and N. Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Transactions on Control of Network Systems*, 5(3):1245–1260, 2018.
- [29] S. S. Ram, A. Nedić, and V. V. Veeravalli. Incremental stochastic subgradient algorithms for convex optimization. *SIAM Journal on Optimization*, 20(2):691–717, 2009.
- [30] P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2012.
- [31] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin. On the linear convergence of the ADMM in decentralized consensus optimization. *IEEE Transactions on Signal Processing*, 62(7):1750–1761, 2014.
- [32] T. Sun, Y. Sun, Y. Xu, and W. Yin. Markov chain block coordinate descent. *Computational Optimization and Applications*, 75(1):35–61, 2019.
- [33] S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- [34] Y. Ye, H. Chen, Z. Ma, and M. Xiao. Decentralized consensus optimization based on parallel random walk. *IEEE Communications Letters*, 24(2):391–395, 2020.
- [35] M. Yemini, R. Saha, E. Ozfatura, D. Gunduz, and A. J. Goldsmith. Semi-decentralized federated learning with collaborative relaying. In *2022 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2022.
- [36] X. Zhang, Y. Liu, J. Liu, A. Argyriou, and Y. Han. D2D-assisted federated learning in mobile edge computing networks. In *2021 IEEE Wireless Communications and Networking Conference (WCNC)*, 2021.
- [37] H. Zhao, B. Li, Z. Li, P. Richtárik, and Y. Chi. BEER: Fast $O(1/T)$ Rate for Decentralized Nonconvex Optimization with Communication Compression. *arXiv:2201.13320*, 2022.